

NCBI Taxonomy Database for Prokaryotic Curation

Shobha Sharma and Conrad Schoch

National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland

Introduction

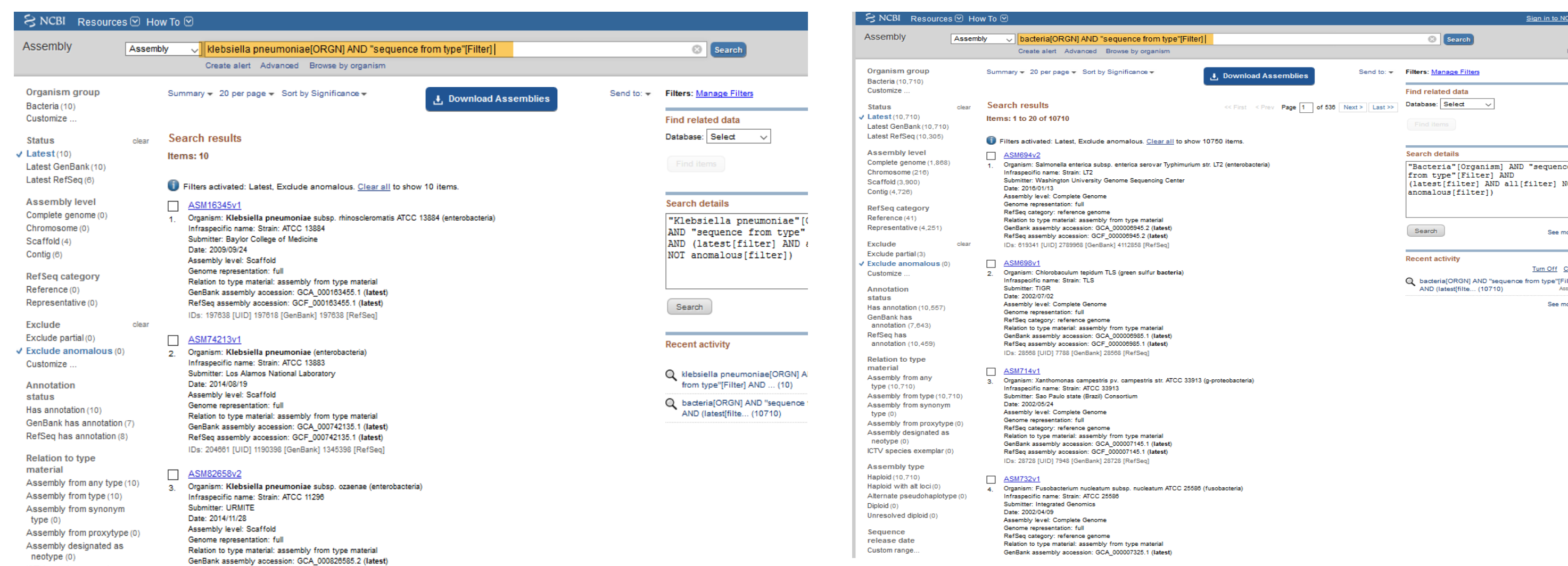
National Center for Biotechnology Information (NCBI) taxonomy database contains the names and phylogenetic lineages of more than **470,000** species with formal names that have molecular data in INSDC databases, about **19,000** of these are species level prokaryote names. Taxonomy database can be regarded as the central organizing hub for many of the resources at NCBI. Here are some of features of NCBI taxonomy database which may not be quite obvious to users.

Type Material in the NCBI Taxonomy Database

Since 2013, GenBank curates type material (including synonym types) in the Taxonomy Database and uses it to flag sequences from type or synonym type in sequence records. Sequence from type is an important subset of GenBank.

Currently GenBank has over **10,700** prokaryotic genome assemblies from type strains.

You can search assemblies from type by using **"sequence from type"[Filter]**

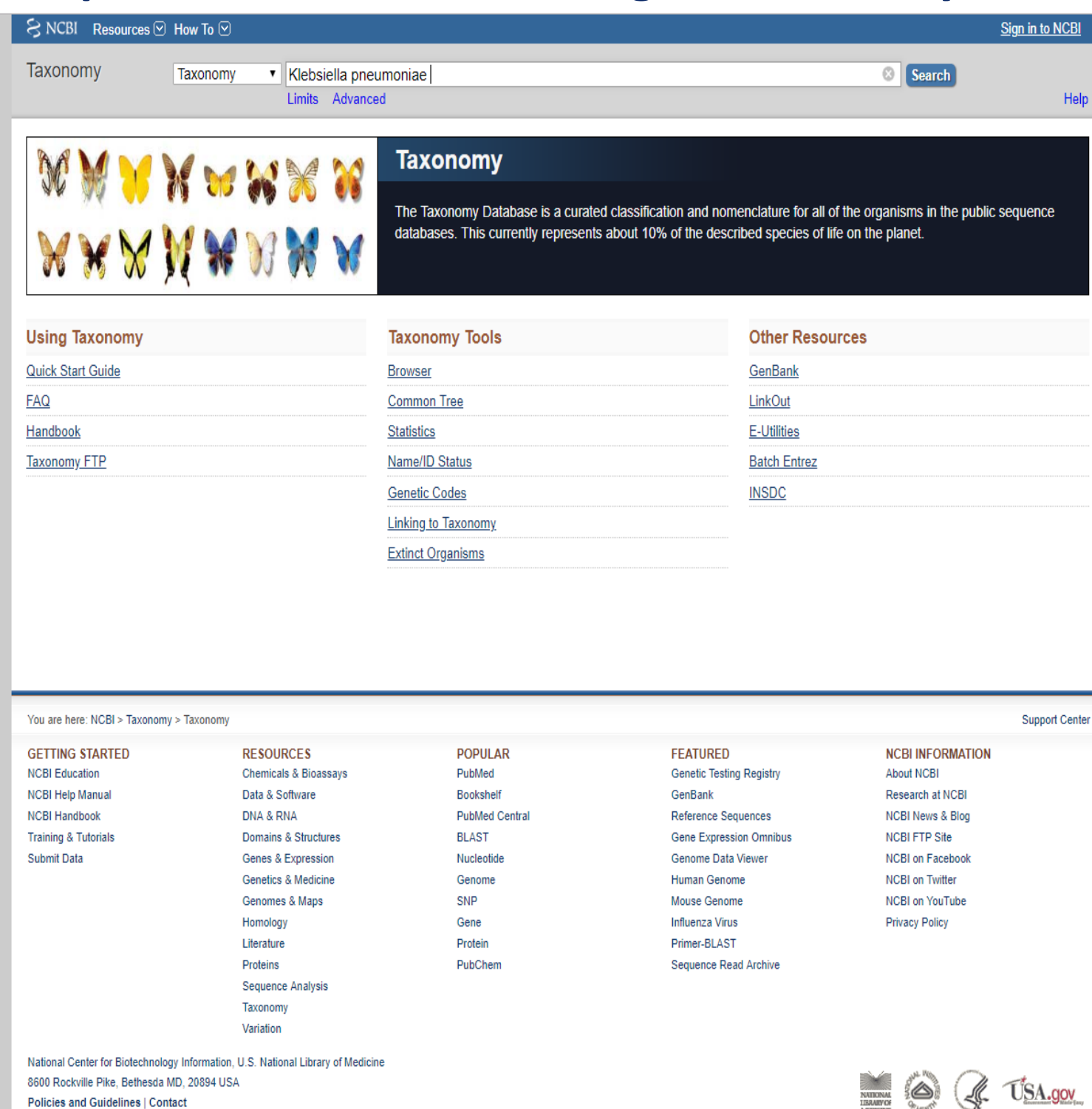


ANI Current Status

- ~10,000 type strains available for ~19,000 species
 - More types added as they are submitted
- ~800 species IDs have been corrected (public assemblies)
- ~2100 new submissions were corrected before they become public
- ~4,700 genome assemblies are currently misidentified
- ANI process executed daily
- Incorporated into PGAP for identification of problems before submission

NCBI Taxonomy Database

NCBI Taxonomy database is the standard nomenclature and classification repository for the International Nucleotide Sequence Database Collaboration (INSDC) which comprises GenBank, ENA and DDBJ databases. It is an entry point into the NCBI system for users who want to find all available information about a particular taxon.



Links to other NCBI databases

Latin words		
Definition	Interim	Decline
Verbale	2100.73	112.87
Pura	10.77 (0.001)	1.03 (0.0001)
Fructus	4.69	2.19
Crucis	2.6	2.19
Papa	2.6	2.19
Caerere/Caerere	2.6	2.19
OS/Occidere	2.6	2.19
Modus/Modus	2.6	2.19
Grat	2.6	2.19
Stil Ligationis	2.6	2.19
Proter/Cutere	2.6	2.19
Indicere/Bona/Gene	1.73 (0.07)	1.68 (0.0001)
Intercept	2.00	2.00
Dirigibile	2.12 (0.0001)	2.00 (0.0001)
Per Verba	2.12 (0.0001)	2.00 (0.0001)
Assensu	2.00	2.00
Papa	2.00	2.00
Per Verba/Intercept	2.12 (0.0001)	2.00 (0.0001)
Intercept	2.00	2.00

Correcting Prokaryotic Genomes Based on Average Nucleotide Identity (ANI)

All new prokaryote genomes submissions to GenBank go through average nucleotide identity (ANI) check to ensure asserted organism name is correct. ANI is also used to correct misidentified genomes that are already public in GenBank

There are **~280,000** live bacterial assemblies in GenBank and **~10,000** assemblies from type stains

Submitted organism: *Vibrio rotiferianus*
 Predicted organism: *Vibrio parahaemolyticus*
 Submitted organism has type: Yes
 Status: MISASSIGNED
 Confidence: HIGH

Accession (NCBI/GenBank)	Organism (Accession)
AN_439 (82.50 Sv Cov)	<i>Vibrio parahaemolyticus</i> (GCA_010585492.2)
NR_038 (82.50 Sv Cov)	<i>Vibrio parahaemolyticus</i> (GCA_01010115.1)
NR_438 (82.49 Sv)	<i>Vibrio parahaemolyticus</i> (GCA_01009875.1)
NR_436 (82.30 Sv)	<i>Vibrio parahaemolyticus</i> NBRC 12711 (GCA_00081305.1)
NR_438 (81.79 Sv)	<i>Vibrio anguillarum</i> (GCA_00002848.2)
NR_853 (69.76 Sv)	<i>Vibrio anguillarum</i> (GCA_00002848.2)
NR_055 (68.70 Sv)	<i>Vibrio anguillarum</i> (GCA_00002848.2)
NR_445 (60.70 Sv)	<i>Vibrio anguillarum</i> (GCA_00017275.1)
NR_446 (60.70 Sv)	<i>Vibrio anguillarum</i> (GCA_00017275.1)
NR_445 (60.70 Sv)	<i>Vibrio anguillarum</i> NBRC 15630 = ATCC 17749 (GCA_00004571.1)
NR_446 (60.70 Sv)	<i>Vibrio anguillarum</i> NBRC 15630 = ATCC 17749 (GCA_00004571.1)
NR_859 (60.39 Sv)	<i>Vibrio anguillarum</i> (GCA_00017275.1)
NR_820 (53.30 Sv)	<i>Vibrio harveyi</i> (GCA_00152585.2)
NR_428 (52.8 Sv)	<i>Vibrio harveyi</i> (GCA_00147157.2)
NR_820 (52.1 Sv)	<i>Vibrio parahaemolyticus</i> (GCA_010585492.2)
NR_820 (52.1 Sv)	<i>Vibrio parahaemolyticus</i> NBRC 15630 = ATCC 17749 (GCA_00004571.1)
NR_820 (51.49 Sv)	<i>Vibrio parahaemolyticus</i> NBRC 15630 = ATCC 17749 (GCA_00004571.1)

Correcting new submissions

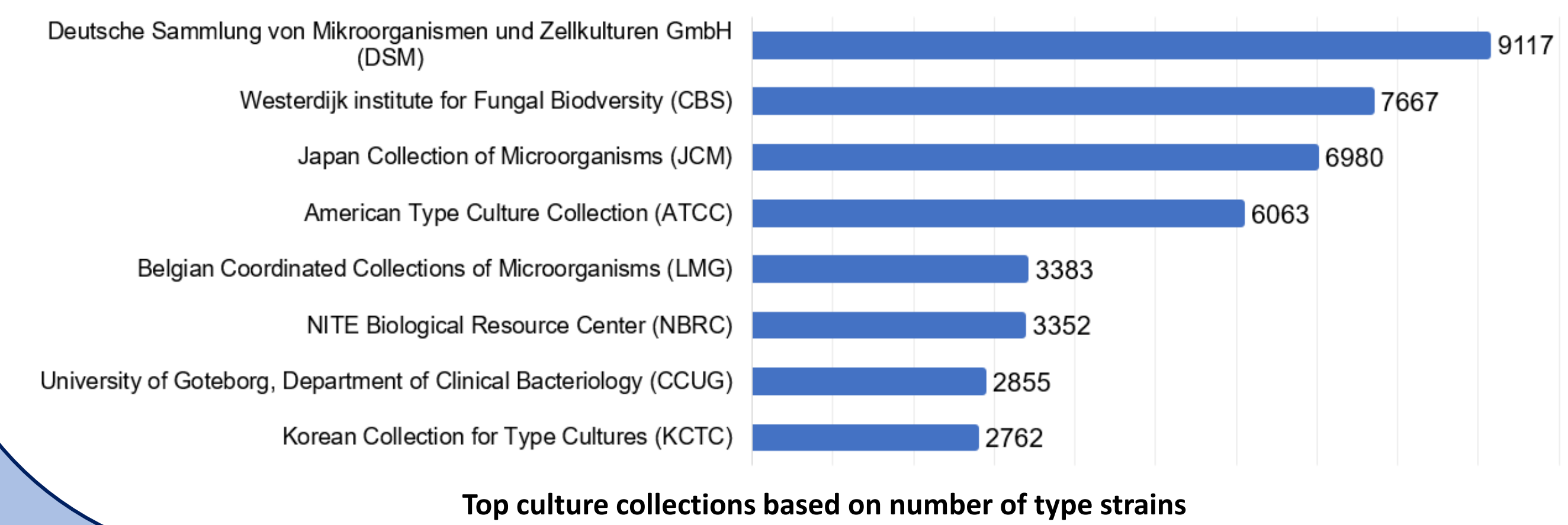
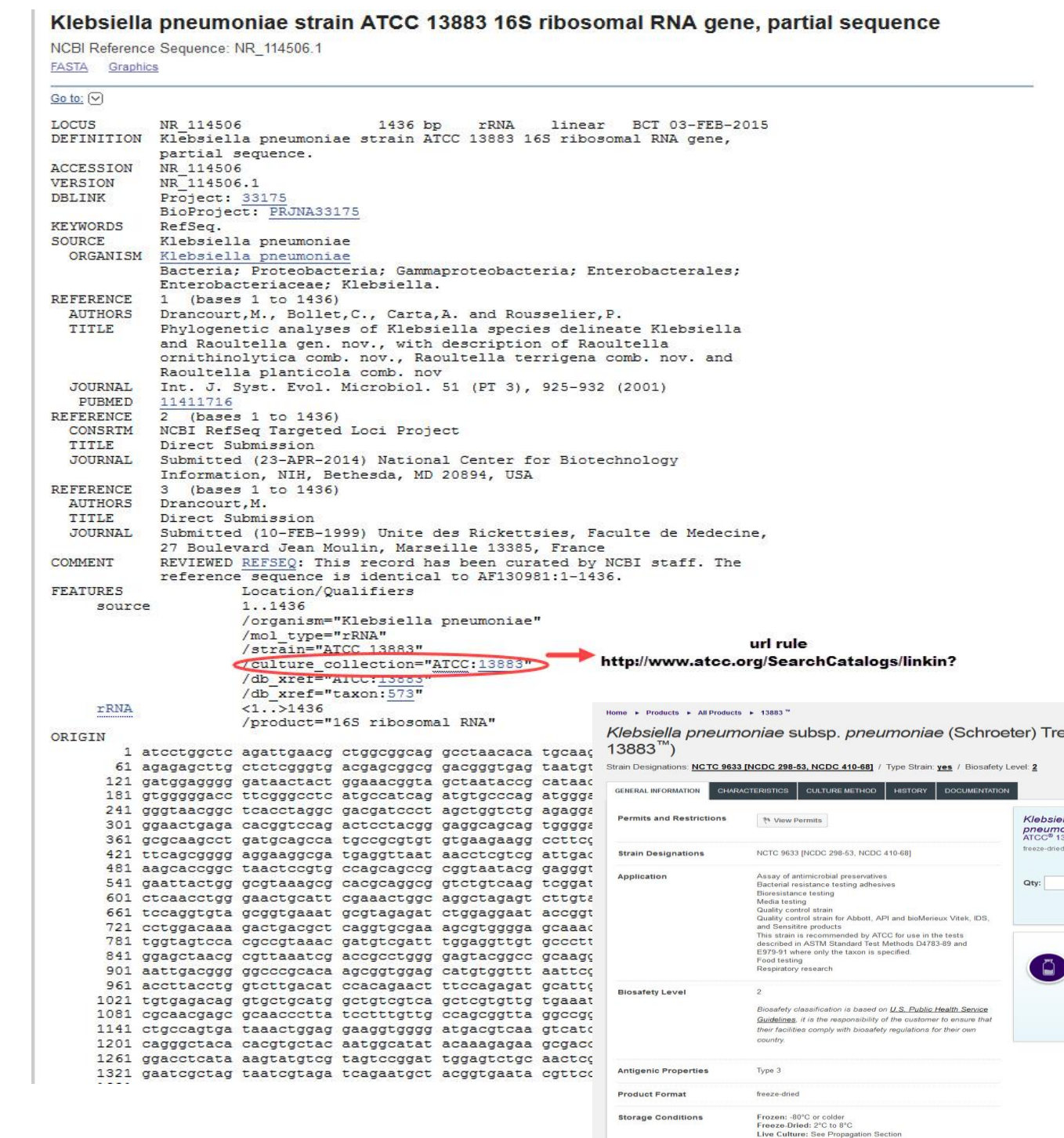
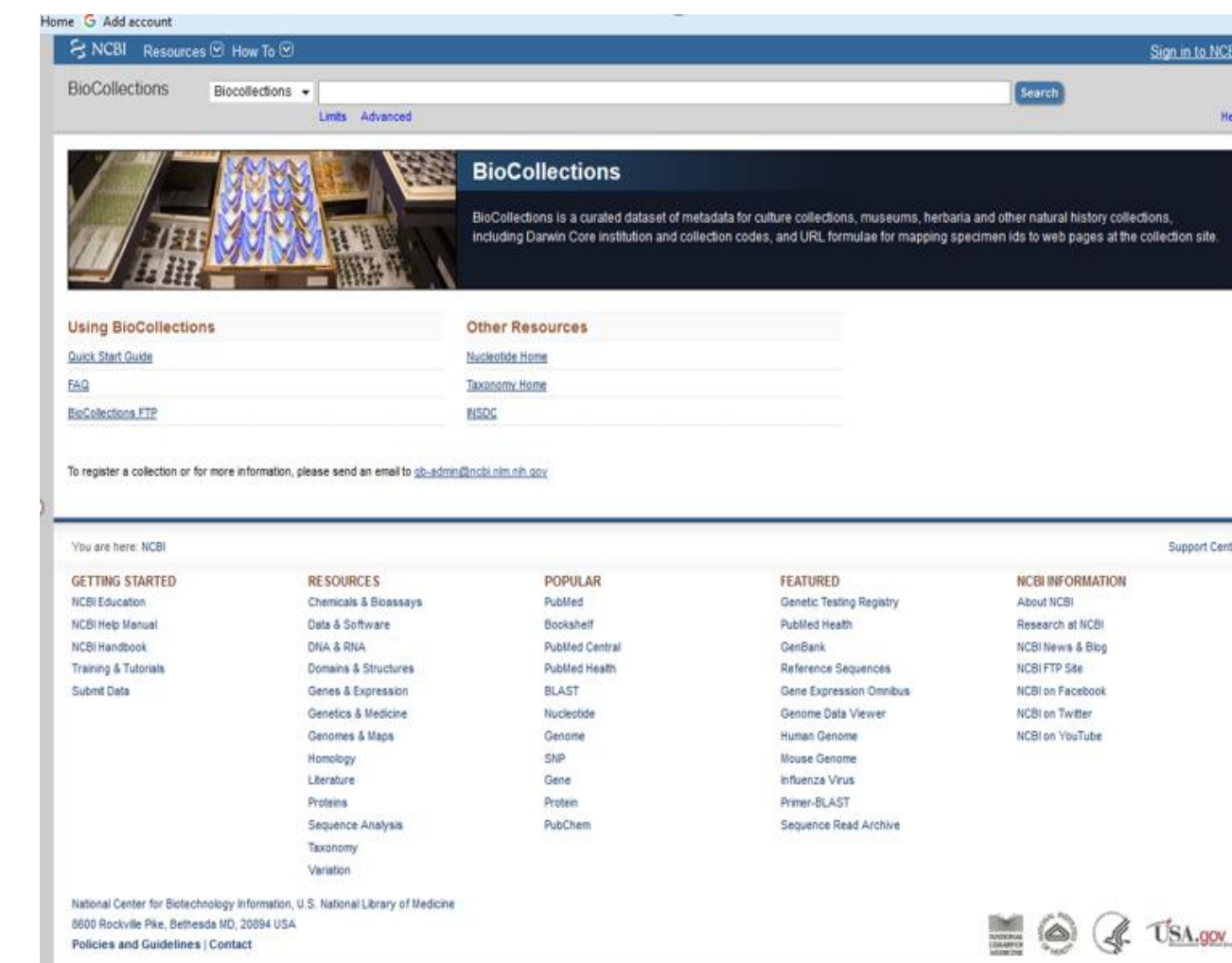
[illegible]

```
##Taxonomic-Update-Statistics-START##
This Genome (Query): 1 GCA_00139855.1
Current Genome (Subject): 1 A015166.1 Arabidopsis thaliana
Previous Name: 1 BR18Arabidopsis longum
Analysis Type: 1 1712
Analysis Type: 1 Average Nucleotide Identity (ANI)
Analysis Type: 1 1712 TITE genome for current name
A1 Genome (subject): 1 GCA_000420505.1
A1 Subject Coverage: 1 98.55561%
A1 ANI: 1 98.55561%
A1 Query Coverage: 1 90%
A1 Subject Coverage: 1 90%
A2 Genome (subject): 1 GCA_001916555.1
A2 Subject Coverage: 1 90%
A2 ANI: 1 98.33813%
A2 Query Coverage: 1 0%
A2 Subject Coverage: 1 0%
##Taxonomic-Update-Statistics-END##
```

NCBI BioCollections Database

Curated dataset of metadata for culture collections, museums, herbaria and other natural history collections connected to sequence records in GenBank. Used to Support the “structured voucher” annotation in the sequence entries submitted to INSDC.

<https://www.ncbi.nlm.nih.gov/biollections>



References

Federhen S. 2012. **The NCBI Taxonomy database**. Nucleic Acids Res. 40: D136-43

Federhen S. 2015. **Type material in the NCBI Taxonomy Database**. Nucleic Acids Res. 43: D1086-98.

Sharma S., Ciuffo S., Starchenko E., Darji D., Chumsky L., Karsch-Mizrachi I. and Schoch CL. 2018. **The NCBI Biocollections Database**. Database(Oxford) doi: 10.1093/database/baz057.

Ciufo S., Kannan S., Sharma S., Badretdin A., Clark K., Turner S., Brover S., Schoch C.L., Kimchi A. and DiCuccio M. 2018. **Using average nucleotide identity to improve taxonomic assignments in prokaryotic genomes at the NCBI.** *Int J Syst Evol Microbiol.* 66:2386-2392. doi: 10.1099/ijsem.0.002809.